
Datasheets for Datasets

1 Links

We attach our web page for the NeurIPS 2021 Datasets and Benchmarks Track reviewers: **Dataset:** <https://bit.ly/38pr1V0>

Our dataset with documentation and associated code will be maintained under the following license: CC-BY-4.0 (<https://creativecommons.org/licenses/by/4.0/>).

2 Author Statement

Our dataset can be commonly copied and redistributed in any medium or format. Our dataset can be accessed using the following license: CC-BY-4.0 (<https://creativecommons.org/licenses/by/4.0/>).

3 Motivation

For what purpose was the dataset created? Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.

Deepfakes have become a challenging well-known technical, social, and ethical issue now. Therefore, many ethical, security, and privacy concerns arise due to the ease of generating deepfakes. And, there is a strong need for developing good deepfake detection methods. However, a good deepfake dataset is required. Although there exist deepfake datasets, there exist no multimodal (video+audio) deepfake dataset. In this work, we aim to fill the gap by providing multimodal (video+audio) deepfake dataset.

Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?

Every author in the paper has contributed.

Who funded the creation of the dataset? If there is an associated grant, please provide the name of the grantor and the grant name and number.

This work was supported by the Basic Science Research Program through National Research Foundation of Korea (NRF) grant funded by Korea government MSIT (No. 2020R1C1C1006004). Also, this research was partly supported by IITP grant funded by the Korea government MSIT (No. 2021-0-00017, Original Technology Development of Artificial Intelligence Industry).

4 Composition

What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people and interactions between them; nodes and edges)? How many instances are there in total (of each type, if appropriate)?

33 Our FakeAVCeleb consists of 20,000+ deepfake videos of 490 celebrities with different ethnic
34 backgrounds belonging to diverse age groups. The real videos that are used to create deepfake dataset
35 are taken from VoxCeleb2 dataset, which consists of real YouTube videos of 6,112 celebrities. . Each
36 ethnic group have a separate directory, and videos are sorted with respect to gender. There are four
37 categories of datasets, for real and three fake, each contains a separate directory.

38 **Does the dataset contain all possible instances or is it a sample (not necessarily random) of**
39 **instances from a larger set?** We include all possible instances or combinations to the best of our
40 knowledge with respect to Audio-Video pairs $((\mathcal{A}_{\mathcal{R}} \mathcal{V}_{\mathcal{R}}), (\mathcal{A}_{\mathcal{R}} \mathcal{V}_{\mathcal{F}}), (\mathcal{A}_{\mathcal{F}} \mathcal{V}_{\mathcal{R}}), (\mathcal{A}_{\mathcal{F}} \mathcal{V}_{\mathcal{F}}))$.

41 **What data does each instance consist of? “Raw” data (e.g., unprocessed text or images) or**
42 **features? In either case, please provide a description.**

43 The dataset consists of processed videos and audios.

44 **Is there a label or target associated with each instance?**

45 Yes. The directory structure of our dataset corresponds to the labels.

46 **Is any information missing from individual instances? If so, please provide a description,**
47 **explaining why this information is missing (e.g., because it was unavailable). This does not**
48 **include intentionally removed information, but might include, e.g., redacted text.**

49 No, all of the annotation information has been provided.

50 **Are relationships between individual instances made explicit (e.g., users’ movie ratings, social**
51 **network links)? If so, please describe how these relationships are made explicit.**

52 Yes. The real videos and generated fake videos from are sorted and named explicitly to in-
53 dicate its type. (e.g., REAL-A/MEN/id00220/iPfJcAJQ84Y/00027.mp4, FAKE-FACESWAP-
54 C/MEN/id00220/00002-id00575-iPfJcAJQ84Y-faceswap-32/00027-fake.mp4)

55 **Are there recommended data splits (e.g., training, development/validation, testing)? If so,**
56 **please provide a description of these splits, explaining the rationale behind them.**

57 Yes, we provide *JSON* files representing each data split for reproducibility of the detection perfor-
58 mance.

59 **Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a**
60 **description.** We tried our best to filter out the bad samples and mistakes with several researchers
61 independently cross-checking samples.

62 **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,**
63 **websites, tweets, other datasets)? If it links to or relies on external resources,**

64 Mostly, FakeAVCeleb is self-contained, however, since our real videos are from VoxCeleb2 and
65 contains only 490 videos, researcher can refer to the VoxCeleb2 dataset if more real videos are
66 required.

67 **Does the dataset contain data that might be considered confidential (e.g., data that is protected**
68 **by legal privilege or by doctor patient confidentiality, data that includes the content of individ-**
69 **uals’ non-public communications)?**

70 No, the dataset is generated real Youtube videos of celebrities. And it does not contain data that
71 might be considered confidential.

72 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening,**
73 **or might otherwise cause anxiety?**

74 No, not at all.

75 **Does the dataset relate to people? If not, you may skip the remaining questions in this section.**

76 Yes. FakeAVCeleb related to people.

77 **Does the dataset identify any subpopulations (e.g., by age, gender)? If so, please describe how**
78 **these subpopulations are identified and provide a description of their respective distributions**
79 **within the dataset.**

80 Yes, the dataset is sorted and each file is named with respect to people with the specific ethnicity and
81 gender.

82 **Is it possible to identify individuals (i.e., one or more natural persons), either directly or indi-**
83 **rectly (i.e., in combination with other data) from the dataset?**

84 Since the videos are of celebrities, it is possible to identify the people. However, there are already
85 openly available dataset information from VoxCeleb2 and YouTube.

86 **Does the dataset contain data that might be considered sensitive in any way (e.g., data that re-**
87 **veals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or union**
88 **memberships, or locations; financial or health data; biometric or genetic data; forms of gov-**
89 **ernment identification, such as social security numbers; criminal history)?**

90 No. The dataset contains YouTube videos of celebrities with different ethnic backgrounds. The
91 intention to include videos of people with diverse ethnic background is to remove racial bias issues in
92 dataset.

93 5 Collection

94 **How was the data associated with each instance acquired? Was the data directly observable**
95 **(e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly in-**
96 **ferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or lan-**
97 **guage)? If data was reported by subjects or indirectly inferred/derived from other data, was**
98 **the data validated/verified? If so, please describe how.**

99 To create deepfake videos, we selected open sourced real YouTube videos from VoxCeleb2 dataset,
100 which consists of real YouTube videos of 6,112 celebrities.

101 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or**
102 **sensor, manual human curation, software program, software API)? How were these mecha-**
103 **nisms or procedures validated?**

104 Three different researchers independently cross-checked the samples and we also used Face++ API
105 to match the best real and fake pairs.

106 **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,**
107 **probabilistic with specific sampling probabilities)?**

108 No, not applicable.

109 **Who was involved in the data collection process (e.g., students, crowdworkers, contractors)**
110 **and how were they compensated (e.g., how much were crowdworkers paid)?**

111 We used the dataset from VoxCeleb2, and three researchers worked on choosing the best samples
112 according to different genders, races, etc, and independently verified the selections one another.

113 **Over what timeframe was the data collected? Does this timeframe match the creation time-**
114 **frame of the data associated with the instances (e.g., recent crawl of old news articles)? If not,**
115 **please describe the timeframe in which the data associated with the instances was created.**

116 The dataset was created from Jan. to June 2021.

117 **Were any ethical review processes conducted (e.g., by an institutional review board)? If so,**
118 **please provide a description of these review processes, including the outcomes, as well as a link**
119 **or other access point to any supporting documentation.**

120 No ethical review process was needed, because all datasets are already available on the Internet.

121 **Does the dataset relate to people? If not, you may skip the remainder of the questions in this**
122 **section.**

123 Yes. It is related to people.

124 **Did you collect the data from the individuals in question directly, or obtain it via third parties**
125 **or other sources (e.g., websites)?**

126 We obtained from VoxCeleb2, online data source.

127 **Were the individuals in question notified about the data collection? If so, please describe (or**
128 **show with screenshots or other information) how notice was provided, and provide a link or**
129 **other access point to, or otherwise reproduce, the exact language of the notification itself.**

130 NA.

131 **Did the individuals in question consent to the collection and use of their data? If so, please**
132 **describe (or show with screenshots or other information) how consent was requested and pro-**
133 **vided, and provide a link or other access point to, or otherwise reproduce, the exact language**
134 **to which the individuals consented.**

135 NA.

136 **If consent was obtained, were the consenting individuals provided with a mechanism to revoke**
137 **their consent in the future or for certain uses? If so, please provide a description, as well as a**
138 **link or other access point to the mechanism (if appropriate).**

139 NA.

140 **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a**
141 **data protection impact analysis) been conducted? If so, please provide a description of this**
142 **analysis, including the outcomes, as well as a link or other access point to any supporting**
143 **documentation.**

144 NA.

145 **6 Preprocessing, cleaning and labeling**

146 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,**
147 **tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing**
148 **of missing values)? If so, please provide a description. If not, you may skip the remainder of**
149 **the questions in this section.**

150 Preprocessing was performed separately for videos and audios. Since we collected videos from
151 VoxCeleb2 dataset, these videos are already face-centered and cropped. We extract respective frames
152 from each video and store them separately, and then extract audios from the videos and store them in
153 *.wav* format with a sampling rate of 16 kHz. Before inputting audio directly to the model for training,
154 we first compute Mel-Frequency Cepstral Coefficients (MFCC) features by applying a *25ms* Hann
155 window [?] with *10ms* window shifts, followed by a fast Fourier transform (FFT) with 512 points.
156 As a result, we obtain a 2D array of 80 MFCC features ($D = 80$) per audio frame and store the
157 resulting MFCC features as a three channel image, which is then passed to the model as an input to
158 extract speech features so that it learns the difference between real and fake human speeches.

159 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support**
160 **unanticipated future uses)? If so, please provide a link or other access point to the “raw” data.**

161 The raw images are available.

162 **Is the software used to preprocess/clean/label the instances available? If so, please provide a**
163 **link or other access point.**

164 Yes, we use Python packages to process our dataset.

165 **7 Use Cases**

166 **Has the dataset been used for any tasks already? If so, please provide a description.**

167 No.

168 **Is there a repository that links to any or all papers or systems that use the dataset? If so, please**
169 **provide a link or other access point. What (other) tasks could the dataset be used for?**

170 We provide the information in the following URL: [http://doi.org/10.23056/FAKEAVCELEB_](http://doi.org/10.23056/FAKEAVCELEB_DASHLAB)
171 [DASHLAB](http://doi.org/10.23056/FAKEAVCELEB_DASHLAB)

172 **Is there anything about the composition of the dataset or the way it was collected and prepro-**
173 **cessed/cleaned/labeled that might impact future uses?**

174 It is explained in our DOI site mentioned above.

175 **Are there tasks for which the dataset should not be used? If so, please provide a description.**

176 No, not applicable.

177 **8 Distribution**

178 **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution,**
179 **organization) on behalf of which the dataset was created? If so, please provide a description.**

180 Use of the dataset is free to all researchers after signing a data use agreement.

181 **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the**
182 **dataset have a digital object identifier (DOI)?**

183 We will provide the dataset through Google drive after receiving the Data Use Agreement (DUA) by
184 Google form. Also, we have DOI: http://doi.org/10.23056/FAKEAVCELEB_DASHLAB

185 **When will the dataset be distributed?**

186 Our dataset is available: http://doi.org/10.23056/FAKEAVCELEB_DASHLAB.

187 **Will the dataset be distributed under a copyright or other intellectual property (IP) license,**
188 **and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU, and**
189 **provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or**
190 **ToU, as well as any fees associated with these restrictions.**

191 Yes. it will be distributed under the following license: Attribution 4.0 International (CC BY 4.0) CC
192 **BY**. This license allows requesters to distribute, remix, adapt, and build upon the material in any
193 medium or format, so long as attribution is given to the creator. The license allows for commercial
194 use.

195 **Have any third parties imposed IP-based or other restrictions on the data associated with the**
196 **instances? If so, please describe these restrictions, and provide a link or other access point to,**
197 **or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these**
198 **restrictions.**

199 NA.

200 **Do any export controls or other regulatory restrictions apply to the dataset or to individual**
201 **instances? If so, please describe these restrictions, and provide a link or other access point to,**
202 **or otherwise reproduce, any supporting documentation.**

203 NA.

204 **9 Maintenance**

205 **Who is supporting/hosting/maintaining the dataset? How can the owner/curator/manager of**
206 **the dataset be contacted (e.g., email address)?**

207 The authors of the paper are maintaining the dataset, particularly by Hasam Khalid:
208 hasam.khalids@g.skku.edu.

209 **Is there an erratum? If so, please provide a link or other access point.**

210 No, not applicable.

211 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete in-**
212 **stances)? If so, please describe how often, by whom, and how updates will be communicated**
213 **to users (e.g., mailing list, GitHub)?**

214 Yes, if there are issues with our dataset, we will immediately reflect them on our dataset site. Also,
215 we keep and provide the updated version with the older version kept around for consistency.

216 **If the dataset relates to people, are there applicable limits on the retention of the data associ-**
217 **ated with the instances (e.g., were individuals in question told that their data would be retained**
218 **for a fixed period of time and then deleted)? If so, please describe these limits and explain how**
219 **they will be enforced.**

220 There are no applicable limits on the retention of the dataset.

221 **Will older versions of the dataset continue to be supported/hosted/maintained? If so, please**
222 **describe how. If not, please describe how its obsolescence will be communicated to users.**

223 Yes, each older version of the dataset will continue to be hosted on FakeAVCeleb dataset site:
224 http://doi.org/10.23056/FAKEAVCELEB_DASHLAB. We add and manage all released versions
225 of our dataset for easy access.

226 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for**
227 **them to do so? If so, please provide a description. Will these contributions be validated/ver-**
228 **ified? If so, please describe how. If not, why not? Is there a process for communicating/dis-**
229 **tributing these contributions to other users? If so, please provide a description.**

230 Others can contact us about incorporating fixes and extensions. We will be happy that our work can
231 be extended further.